# RefSeq Functional Elements: An Annotated Dataset of Validated Non-Genic Elements for Variant Interpretation and Functional Discovery in the Human Genome

Catherine M. Farrell, Terence D. Murphy, and the RefSeq Curation and Development Team*

National Center for Biotechnology Information (NCBI), National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA

Follow NCBI on:

**Abstract:** The human genome contains many non-genic elements that play roles in gene regulation, chromosome organization, recombination, repair or DNA replication. Human disease can result from sequence variation in those elements, with many genome-wide association studies indicating disease-associated variation in non-coding regions. The locations of gene regulatory elements can be predicted from several large-scale epigenomic mapping projects, but those data are not generally visible in traditional genome annotation, are difficult to interpret in the absence of specialized research knowledge or customized displays, and do not always show function when tested experimentally. NCBI has therefore introduced a more accessible dataset, RefSeq Functional Elements (www.ncbi.nlm.nih.gov/refseq/functionalelements/), which are annotated on the human genome alongside conventional genes. This curated dataset, which is restricted to known elements from published experimental data, includes richly annotated RefSeq records and accompanying descriptive records in the Gene database (www.ncbi.nlm.nih.gov/gene/). The dataset includes known enhancers, silencers, recombination regions, and other non-genic regions with experimentally-validated function. As of NCBI's Updated Annotation Release 109.20190905 on the GRCh38.p13 genome assembly, the dataset includes over 3.8K GeneIDs and 8.9K feature annotations, with further growth expected for future NCBI annotation releases. The dataset is publicly available for FTP download (ftp://ftp.ncbi.nlm.nih.gov/genomes/Homo_sapiens/). Feature annotation can be visualized in the 'Biological regions' track available in NCBI browsers, including the Genome Data Viewer. This data track is particularly useful when viewed alongside other available tracks, such as variation tracks from dbSNP, ClinVar or dbVar, or study-specific custom tracks or track hubs. These non-genic annotations provide insights into non-coding genome function. They are valuable for basic discovery of gene regulatory regions, interpretation of non-coding variants, or as known positive controls for genome-wide studies aimed at discovering additional elements.

## 1. RefSeq Functional Elements definition and scope

**Definition:**
- Any **non-genic** genomic element that has functional significance **based on experimental support**, and is not otherwise considered a conventional gene
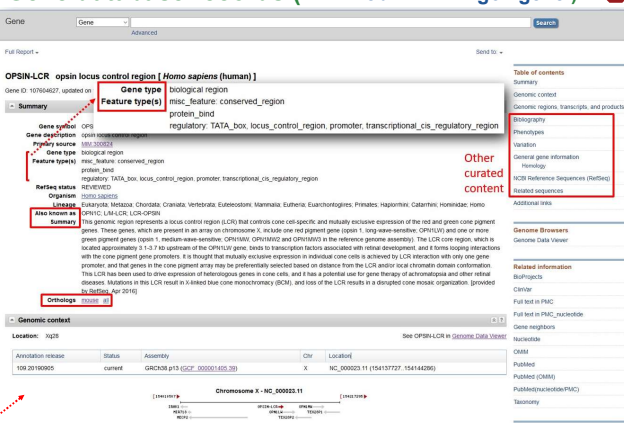
**Types:**
- **Gene regulatory elements**, e.g. enhancers, silencers, promoters, protein binding sites
- **Known structural elements**, e.g. boundary elements, matrix/scaffold-associated regions, other structural regions associated with chromatin conformation
- **Other elements of functional importance**, e.g. clinically-significant recombination hotspots, well-defined replication origins

**Annotation scope:**
- Functional elements that have been **experimentally validated**
- **Human** and **mouse** elements
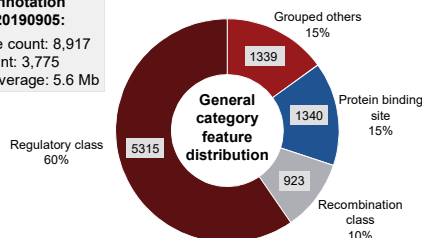- Priority for **elements implicated in human disease**

## 2. Gene database records (www.ncbi.nlm.nih.gov/gene/)



OPSIN-LCR opsin locus control region [ *Homo sapiens* (human) ]

Other curated content

## RefSeq records

- Genomic accessions with **NG_** prefixes
- Associated with NCBI **BioProject accession PRJNA343958**
- International Nucleotide Sequence Database Collaboration (INSDC) feature annotation
- Additional controlled vocabularies from the Sequence Ontology
- Experimental evidence qualifiers with:
  - the type of **evidence** (Evidence & Conclusion Ontology codes)
  - links to **publications** in PubMed
- Additional qualifiers with further details from the literature, including /function, /note or /bound_moiety qualifiers

```
regulatory    867..22391
/regulatory_class="locus_control_region"
/experiment="EXISTENCE:transgenic organism evidence
[ECO:0001131][PMID:3690667]"
/note="21.5 kb ClaI-BglII fragment from -1 kb to -22.5 kb
includes 5'HS1-5'HS5"
/function="regulates developmental expression of the
beta-globin genes"
/db_xref="GeneID:109580095"
regulatory    3114..5057
/regulatory_class="matrix_attachment_region"
/experiment="EXISTENCE:fractionation evidence
[ECO:0000100][PMID:3208739, PMID:2559410]"
/note="2 kb AvaII pGSE228 fragment c SAR"
/function="associates with the nuclear scaffold"
/db_xref="GeneID:109580095"
regulatory    5274..6209
/regulatory_class="enhancer"
/experiment="EXISTENCE:reporter gene assay evidence
[ECO:0000049][PMID:20231293, PMID:17548470]"
/note="5'HS1 or HS1-3' enhancer fragment"
/function="enhancer in K562 erythroleukemia cells"
/db_xref="GeneID:109580095"
```

*Example: Segment of NG_052895.1, HBB-LCR, GeneID:109580095*

View RefSeq Functional Elements **feature annotation** in an NCBI genome browser, including the Genome Data Viewer (ncbi.nlm.nih.gov/genome/gdv/) or the Variation Viewer (ncbi.nlm.nih.gov/variation/view/):

- Configure to display the **'Biological regions'** track (see top right of image)
- Load NCBI 'Genes' or variation tracks (e.g. dbSNP, ClinVar, dbVar), configure a track hub of interest, or use the 'Custom Data' tab in the configuration interface to **view alongside other data of choice**
- Mouse-over a Functional Element feature to view the **associated tooltip** with relevant functional data, including links to publications and sequences

*Example: Genome Data Viewer browser image showing LOC107303343 feature annotation overlapping the ADA gene and variation data on chromosome 20. Tooltips are displayed for two select features.*



## 3. Data access

Find multiple ways to access our data on our website:
ncbi.nlm.nih.gov/refseq/functionalelements/

**NCBI RefSeq Functional Elements**
- Overview
- RefSeq Functional Element Records
  - RefSeq Functional Element Feature Annotation
  - Feature Annotation Glossary
- Data Access
  - Access via Gene
  - Access via Nucleotide
  - Access via BLAST
  - Access via BioProject
  - Access via NCBI Graphical Displays
  - Access via FTP
    - Feature table
    - Feature extraction examples
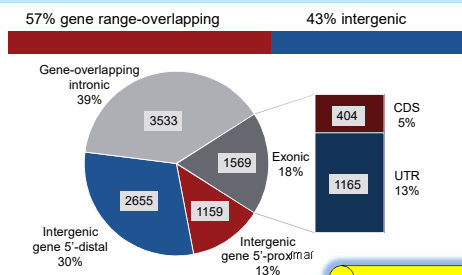
## 5. NCBI genome browser view



## Analysis of RefSeq Functional Elements

**From NCBI Annotation Release 109.20190905:**
- Total feature count: 8,917
- GeneID count: 3,775
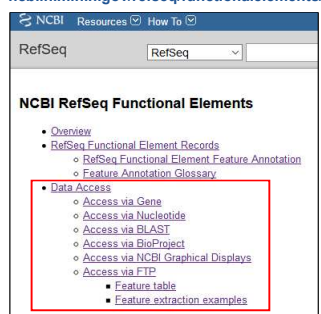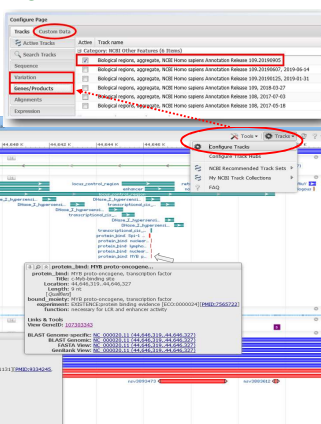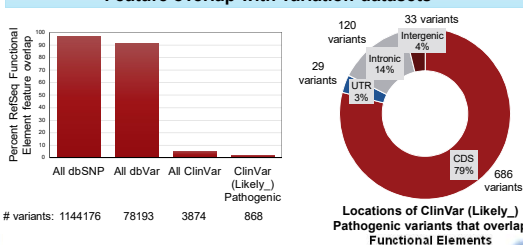- Genome coverage: 5.6 Mb



**General category feature distribution**
- Grouped others 15% — 1339
- Protein binding site 15% — 1340
- Recombination class 10% — 923
- Regulatory class 60% — 5315

### Feature locations relative to genes

57% gene range-overlapping | 43% intergenic



- Gene-overlapping intronic 39% — 3533
- Exonic 18% — 1569
- Intergenic gene 5'-proximal 13% — 1159
- Intergenic gene 5'-distal 30% — 2655
- CDS 5% — 404
- UTR 13% — 1165

**Overlapping gene biotypes:**
- 74% protein-coding
- 22% lncRNA
- 3% pseudogene types
- 1% others

**They include:**
- 797 RefSeqGene (RSG) genes
- 122 Locus Reference Genomic (LRG) genes
- 638 genes used for ClinVar submissions

> - Significant gene overlap is observed
> - Gene regions may not always have exclusive transcript or protein-coding function!

**Gene association caution:**
> - Gene regulatory elements don't necessarily regulate the genes they overlap or are closest to
> - *Example:* the *ZRS* limb enhancer (GeneID:105804841) in an intron of the *LMBR1* gene regulates *SHH* located 1 Mb away!

### Feature overlap with variation datasets



# variants: All dbSNP 1144176 | All dbVar 78193 | All ClinVar 3874 | ClinVar (Likely_) Pathogenic 868

- Intergenic 4% — 33 variants
- Intronic 14% — 120 variants
- UTR 3% — 29 variants
- CDS 79% — 686 variants

**Locations of ClinVar (Likely_) Pathogenic variants that overlap Functional Elements**

**Implications for variant interpretation:**
> - Not all disease-associated variation may be due to alterations in protein-coding or transcript function. Mutations in overlapping elements with non-transcript-related function may be disease-causing too!

Also check out NCBI's **gene/transcript-related RefSeq** annotation, including the new **Matched Annotation from NCBI and EMBL-EBI** (MANE) dataset: www.ncbi.nlm.nih.gov/refseq/MANE/

**www.ncbi.nlm.nih.gov/refseq/functionalelements/**